



Deutsche Bahn

The not-so-efficient train network.

Iron Hack Final Project - Ray Pham

IRON
HACK

About Deutsche Bahn

Deutsche Bahn is the national railway company of Germany.

It is the largest railway operator in Europe and the second-largest in the world by revenue.

Deutsche Bahn is a major employer in Germany with over 324,000 employees.

Generates over €56B in revenue per year.



Quick Facts

It operates a network of over 33,000 kilometres of track. This network includes over 3,000 stations and serves over 2,000 cities and towns.

It operates more than 40,000 trains daily.
Tickets for the long-distance service range
between 29€ to 150€ per journey.

Deutsche Bahn carries almost 2 billion passengers per year.



German Not-so-efficient

Deutsche Bahn once had a reputation akin to the entire image of Germany, punctuality and efficiency.

However in the past years, due to the aging and neglected infrastructure, this has become a problem for train reliability and punctuality.

80,6%

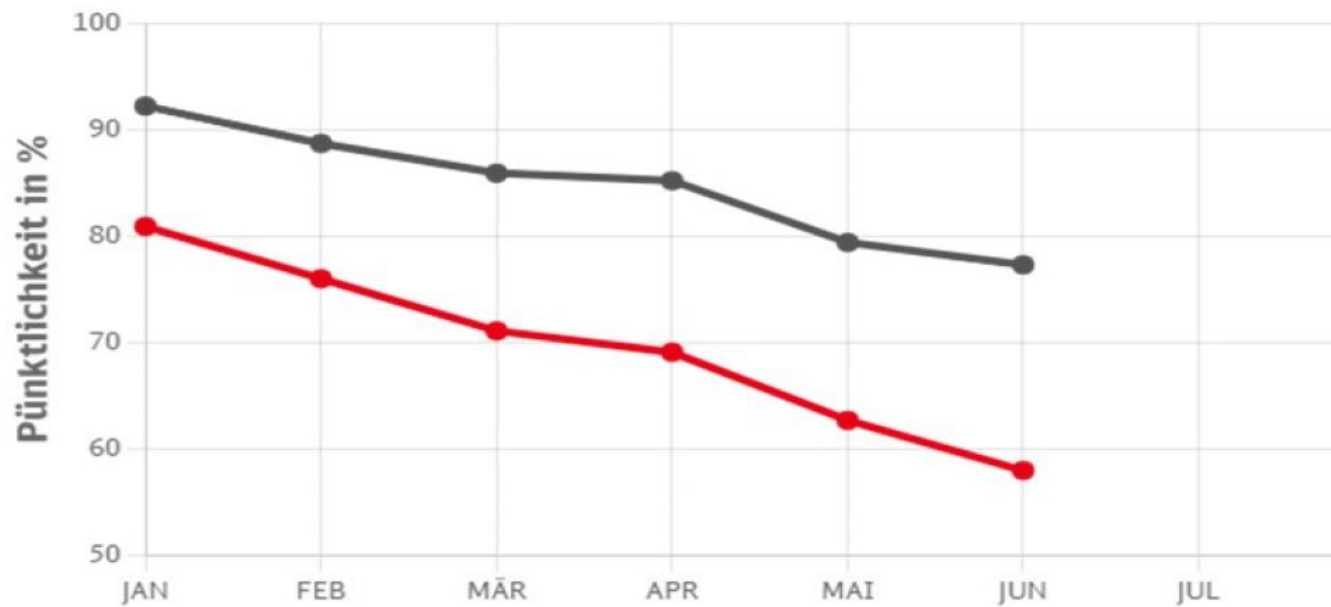
long-distance
trains on time
in 2022.

58%

long-distance
trains on time
in 2023.

In Germany, a train is considered 'on-time' if it arrives no later than 6 minutes of its scheduled time.

2023 January to June



Objectives

This analytic project is to understand the efficiency of the Deutsche Bahn train network on key routes, by analysing key dimensions such as delays, distance, timings.

1. Which stations have the most delays
2. If time of travel causes delays
3. Does distance of between affect delays
4. Build a model that can predict potential delays

Unfortunately I wasn't able to retrieve the passenger information and due to the amount of travels and stations, I focused the project on key cities.

The data observed

Period: May - August 2022

Trips: 428595 (Time/Date)

Routes: 704 (Depart/Arrive Station)

Stations: 36 (Germany and EU)

Distance Calculated (Lat/Lon)

Train Types:

ICE = InterCity Express (DE)

IC = InterCity (DE)

EC = EuroCity (EU)

FLX = FlixTrain (EU)

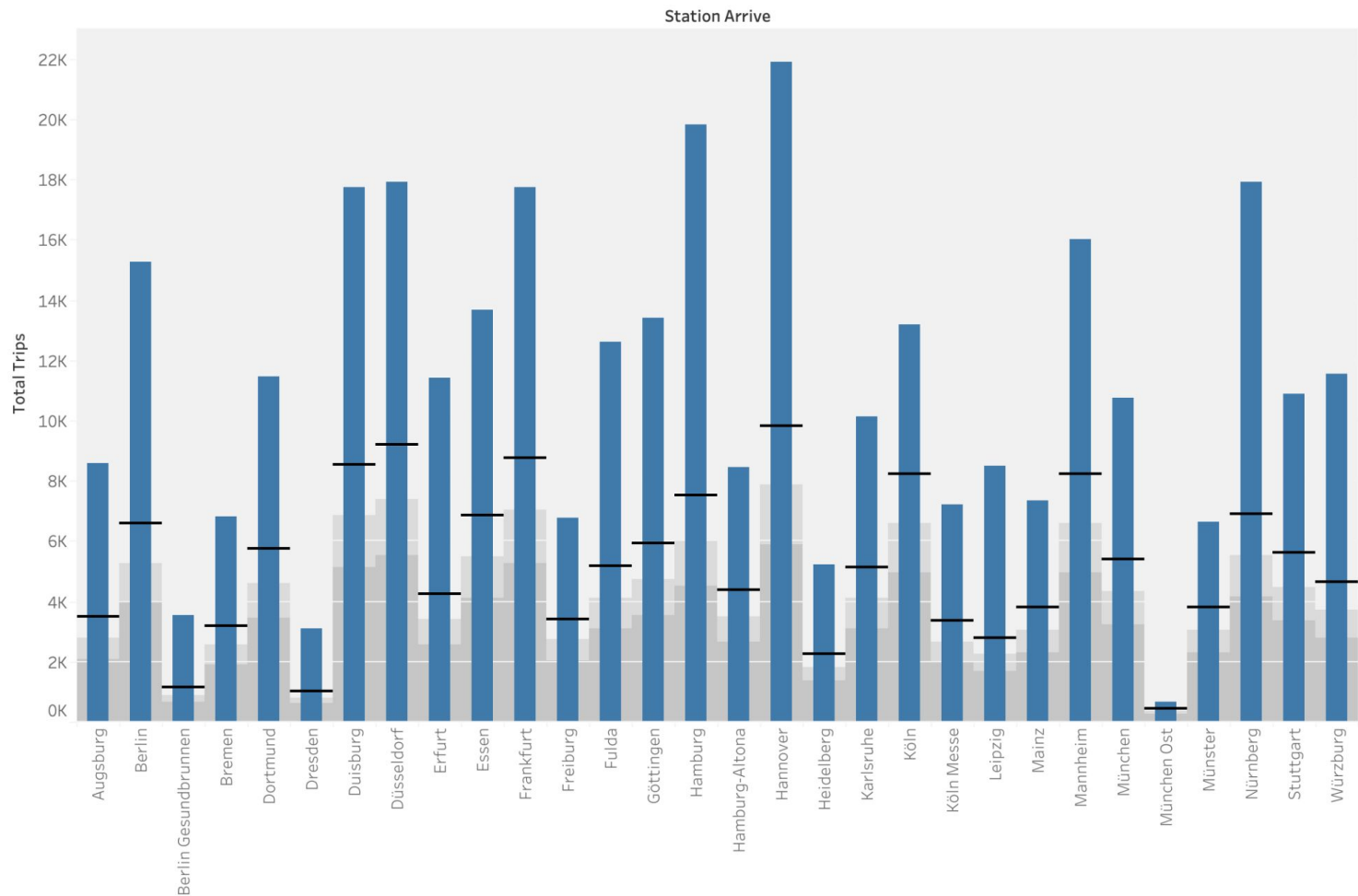
ECE = EuroCity Express (DE)

EN = EuroNight (EU)

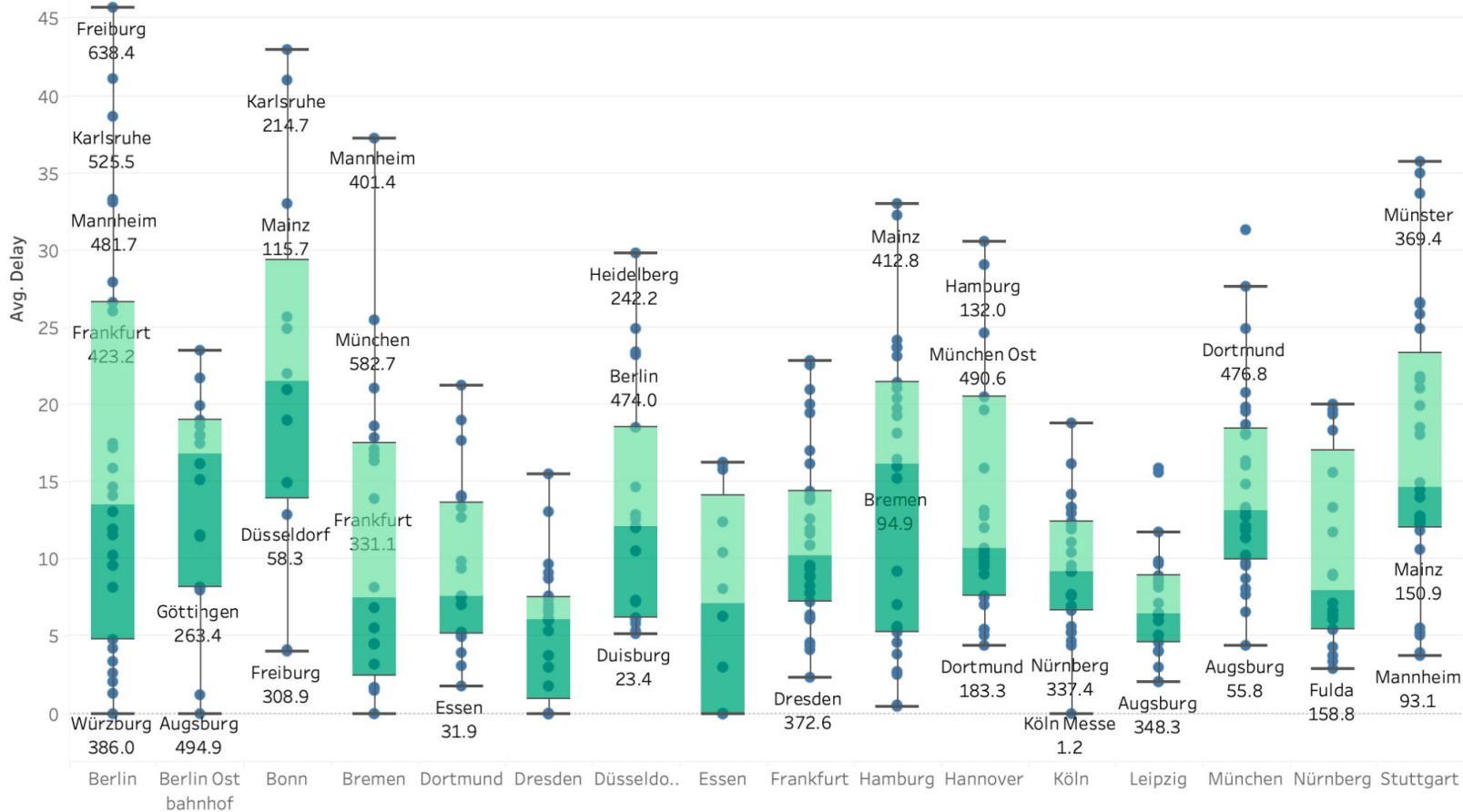


Top 10 delay routes

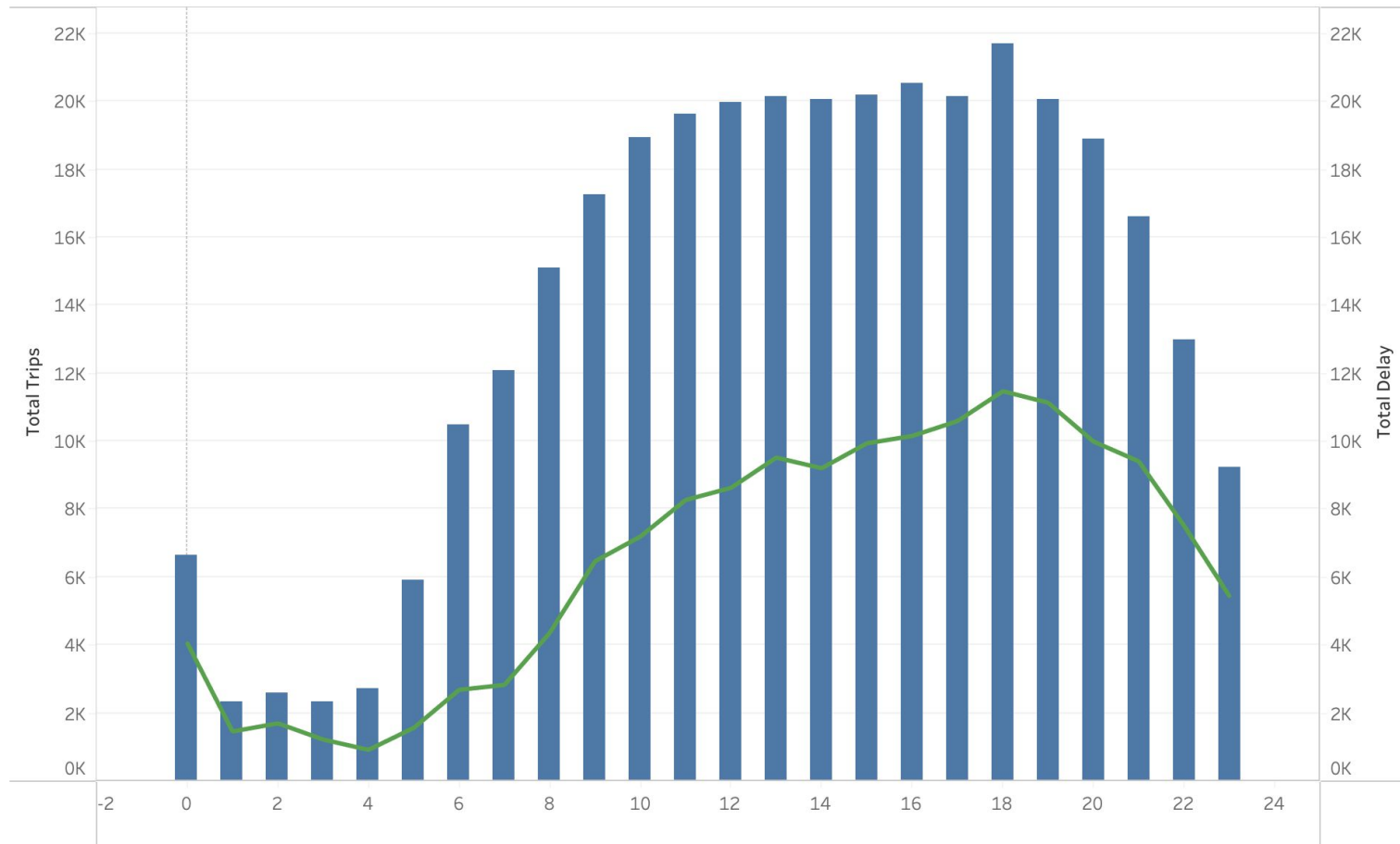
Arrival	Departure	Trips	Delays	Delay %	Avg. Min	Distance KM
Hamburg	München	4826	2649	54.89	28.22	611.578
Nürnberg	München	7357	2321	31.55	18.82	149.285
Hamburg-Altona	München	4204	2263	53.83	26.81	612.556
München	Hamburg-Altona	4268	2237	52.41	29.17	612.556
Mannheim	Hamburg-Altona	2921	2212	75.73	31.8	464.138
Hannover	Köln	3394	2161	63.67	24.67	249.447
Nürnberg	Hamburg-Altona	4198	2037	48.52	25.87	463.588
Hannover	Hamburg-Altona	4116	1983	48.18	19.77	131.495
Düsseldorf	München	2669	1965	73.62	25.93	484.48'
Frankfurt	München	3512	1952	55.58	21.99	303.663



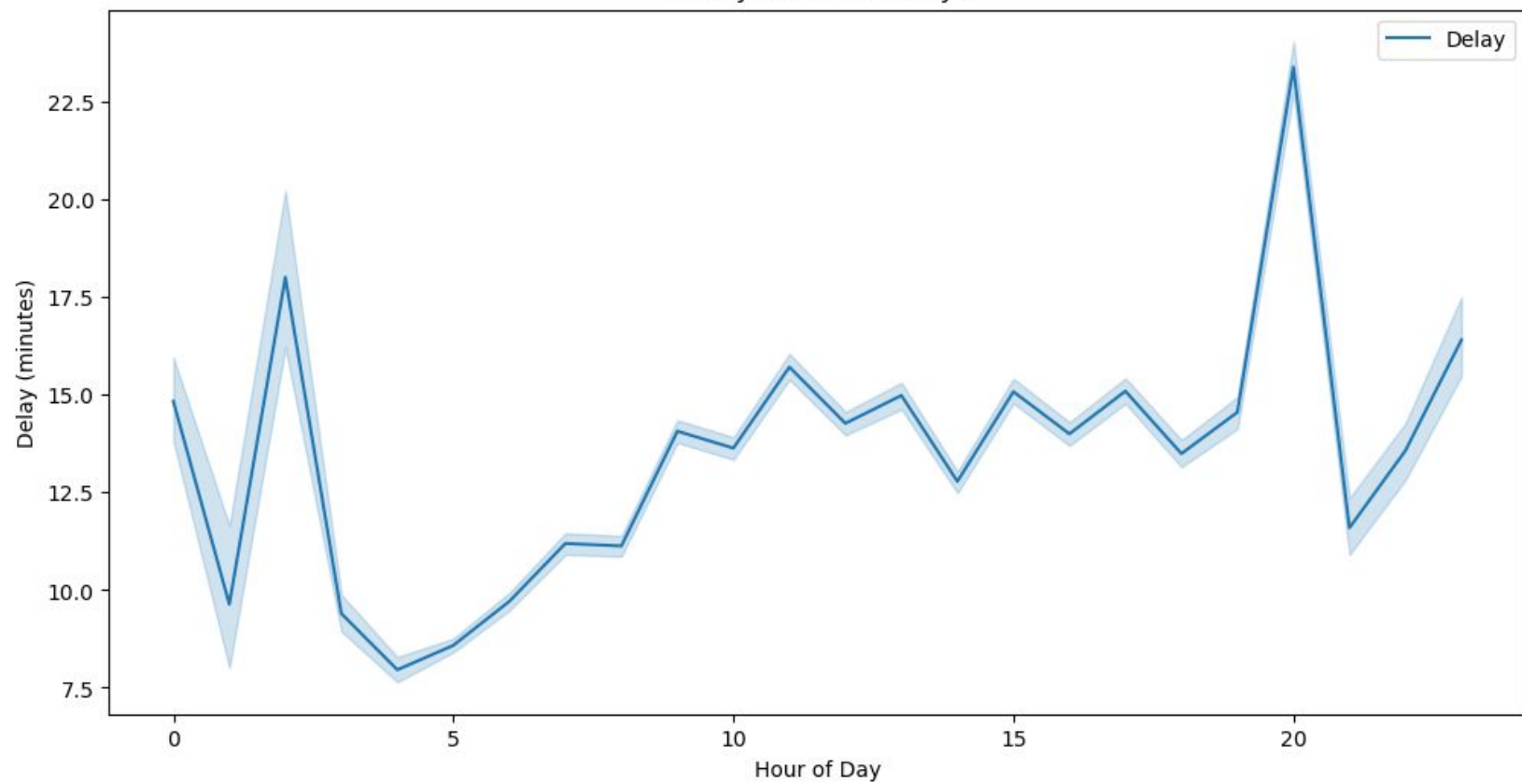
Delays between German cities



Hourly Trends



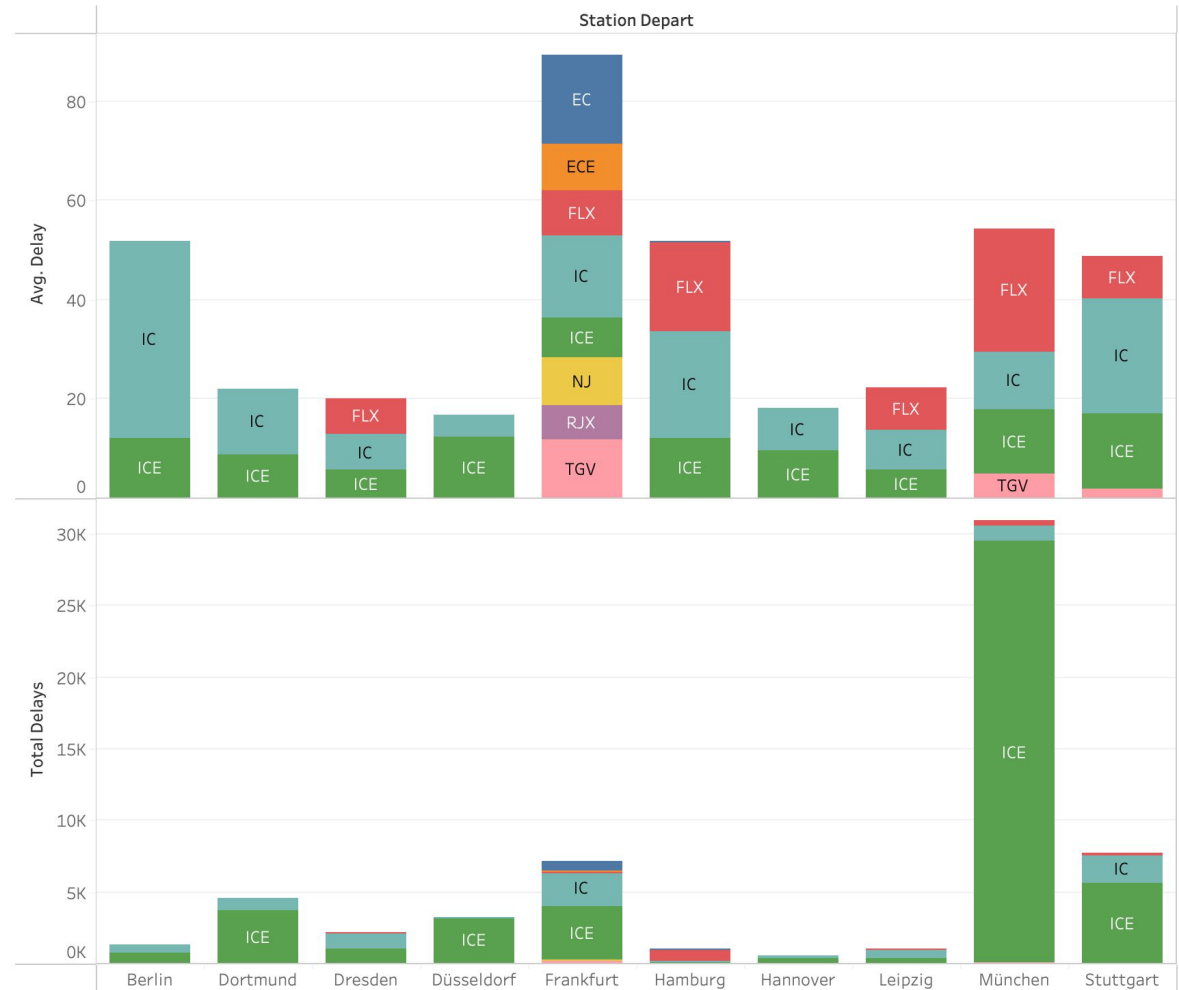
Hourly Trends of Delays



Train Types

Observing the types of trains with the most delays, we can see that the ICE service from Munich contributes to the majority of delays.

But on average, the delay time is roughly equal across all station departures.



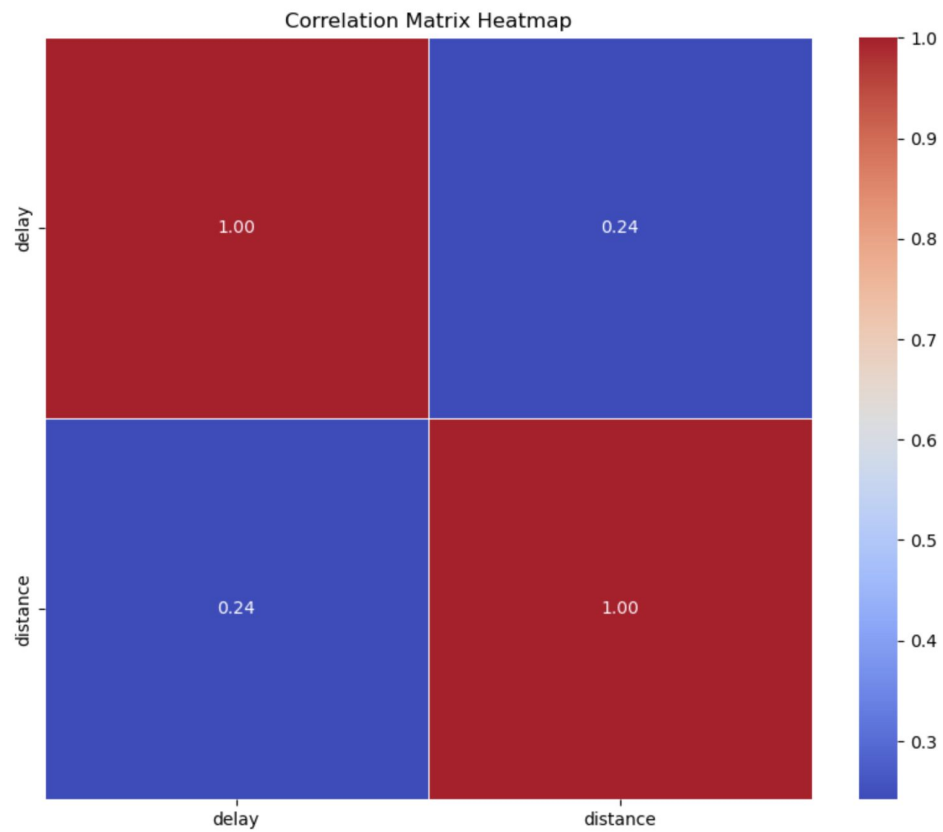
Station Depart	Train Type	Trips	Avg. Delay	Total Delays	Canceled
Berlin	IC	24,286	40	532	6
	ICE	22,698	12	760	143
Dortmund	IC	23,609	13	875	150
	ICE	92,978	9	3,703	402
Dresden	FLX	2,064	7	116	6
	IC	23,933	7	1,013	164
	ICE	22,588	6	1,091	201
Düsseldorf	IC	58	4	5	3
	ICE	85,188	12	3,228	532
Frankfurt	EC	17,486	18	657	84
	ECE	3,437	9	151	72
	FLX	996	9	22	65
	IC	63,821	16	2,311	242
	ICE	89,839	8	3,730	1,033
	NJ	648	10	26	2
	RJX	1,117	7	54	11
Hamburg	TGV	5,365	12	240	13
	EC	7	0	0	0
	FLX	35,227	18	878	127
	IC	2,006	21	50	1
	ICE	3,835	12	141	6
Hannover	IC	6,042	9	194	13
	ICE	9,308	10	389	62
Leipzig	FLX	2,951	8	94	12
	IC	13,749	8	561	86
	ICE	7,366	6	368	60
München	FLX	14,388	25	337	39
	IC	23,750	12	1,035	238
	ICE	825,258	13	29,475	2,734
	TGV	1,582	5	72	8
Stuttgart	FLX	8,184	9	249	30
	IC	67,202	23	1,861	179
	ICE	170,189	15	5,636	1,033
	TGV	13	2	1	0

Model Testing

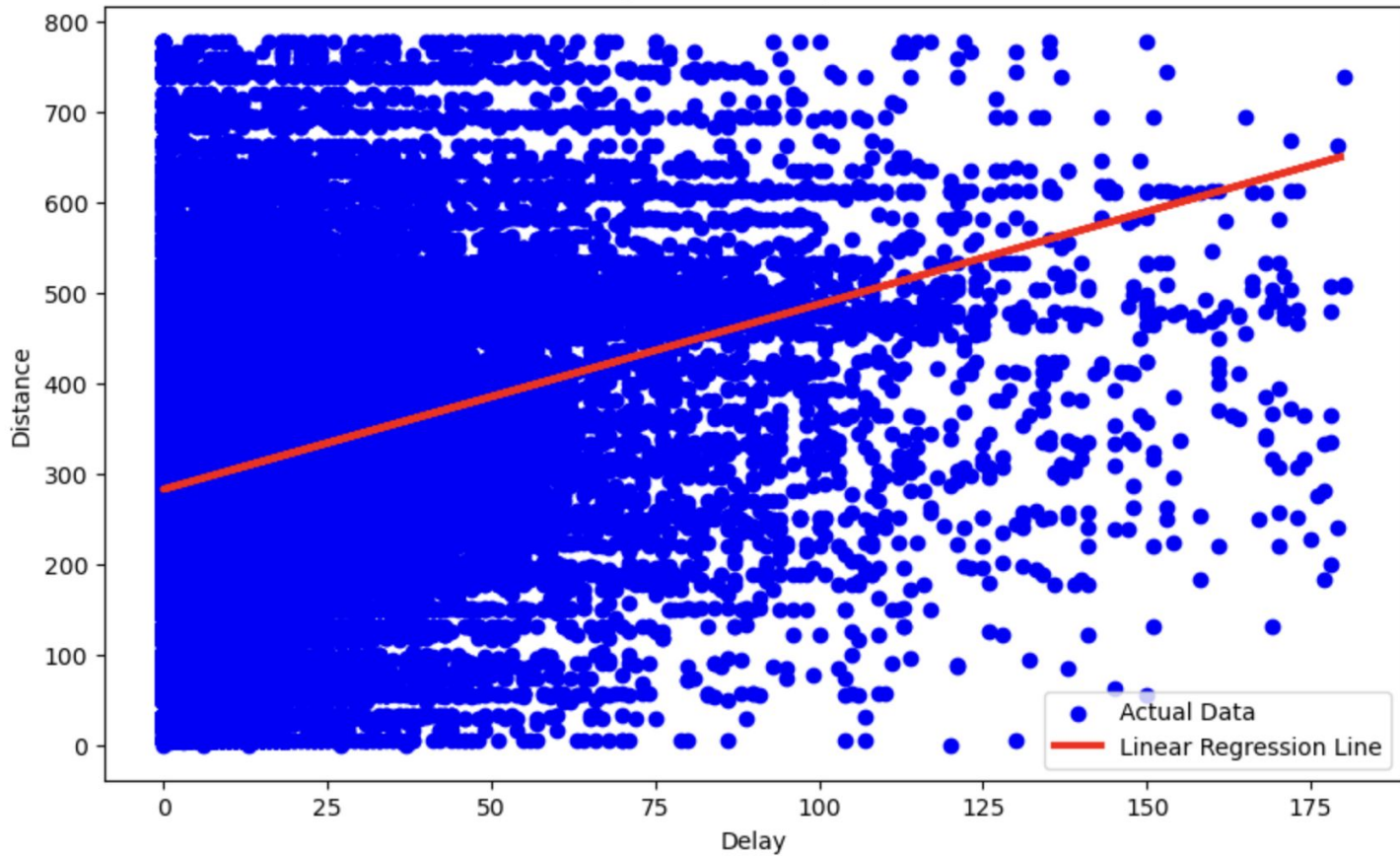
Looking at the initial analysis, I can already observe that there isn't a very clear pattern between features and numericals, so tested a few different models to see, if based on one numerical value (distance) could be enough to build a model.

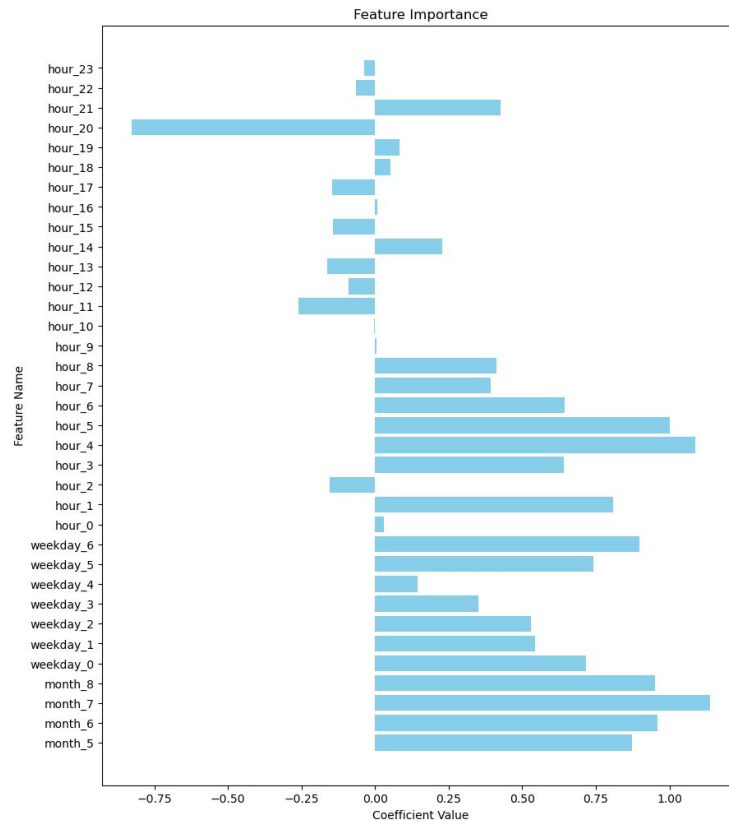
For the test for the features was used

- Station (Destination & Arrival)
- Day of Week
- Hour of Day
- Type of Train (ICE, IC etc)



Linear Regression: Delay vs. Distance





Model Testing

Logistic
MSE:
372.4184

Random Forest
MSE:
464.4855

Decision Tree
MSE:
525.8380

Really bad scores on the model testing.

Testing a Classification Model

	precision	recall	f1-score
0	0.71	0.72	0.71
1	0.67	0.65	0.66

To approach a different model, 0 = No, 1 = Yes binary classification was added if the the percentage if the train will be delayed for longer than 6 minutes.

First Learnings

My first learnings is with the dataset, and after the wrangling i performed, it didn't do enough to build a model that was able to predict delays on an accurate basis.

Things I could try to apply as next steps:

- Enrich the dataset with more variables (such as passenger count per station per day)
- Test different approaches with dates (clustering between weekday and weekend).
- Add an additional numerical dimension by calculating the total number of trips between key certain.

Conclusion

My assumption of distance impacts delays applies in some cases only. Time of travel has a minor effect on the likelihood that a train will be late (early morning as opposed to rush hour which also affects the average delay duration) but the impact follows the trip frequency trend.

More time and data is needed, as well as finding a way to cluster the data that would allow myself to build a better prediction model.

Train Delay Prediction



Under Construction - Just like the Deutsche Bahn



Departure Station:

Arrival Station:

Day of Week:

Hour of Day:

Danke.